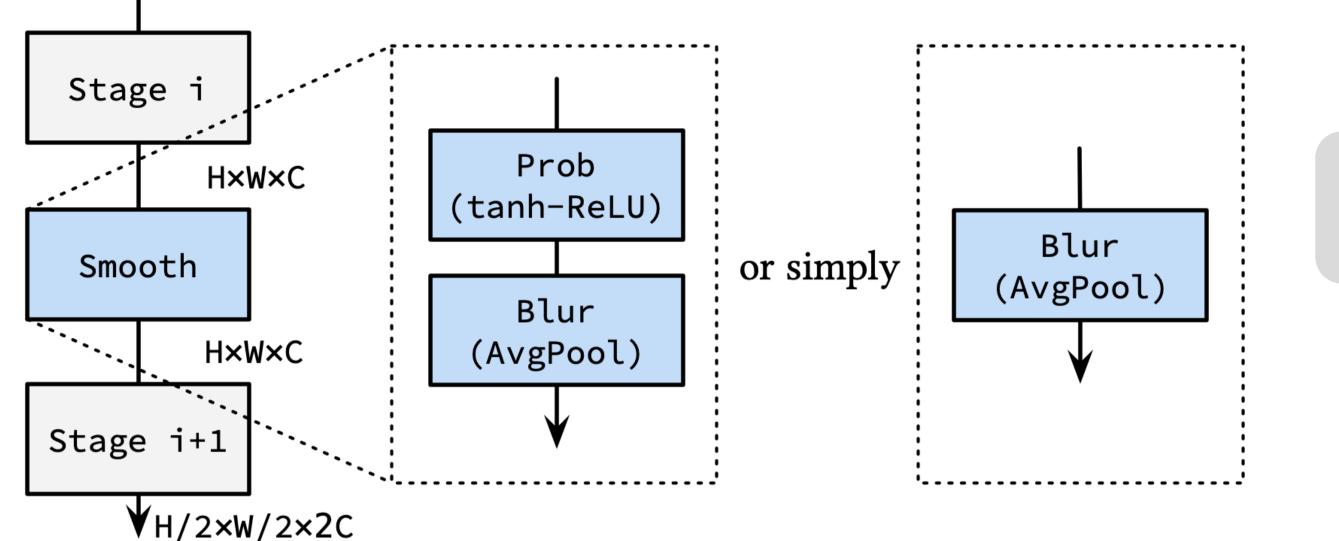
# Blurs Behave Like Ensembles

We introduce a novel ensemble method, "spatial ensemble". Spatial ensemble improve accuracy, uncertainty, and robustness without increasing inference time. But what is a spatial ensemble? Spatial ensemble, or spatial smoothing, is a method that aggregate nearby feature maps (See  $\searrow$ ).

# I. How Can We Apply Spatial Smoothing to NNs?



## Figure 1. Simply add

AvgPool2d(kernel\_size=2,
 stride=1, padding=1)

at the end of stages. For example, use four spatial smoothing layers for ResNet.

# II. Spatial Smoothing Significantly Improves MC Dropout

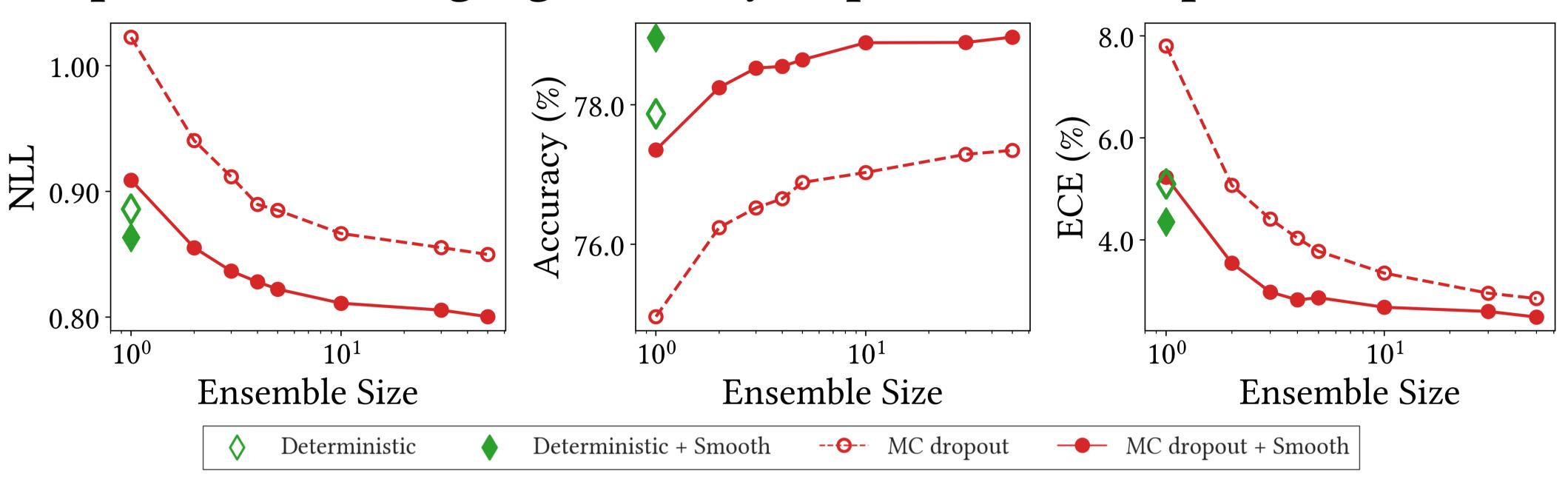


Figure 2. "MC dropout + spatial smoothing" is **25**× **faster** than canonical MC dropout with similar predictive performance. Spatial smoothing also improves deterministic NNs.

# III. How Does Spatial Smoothing Improve NNs?

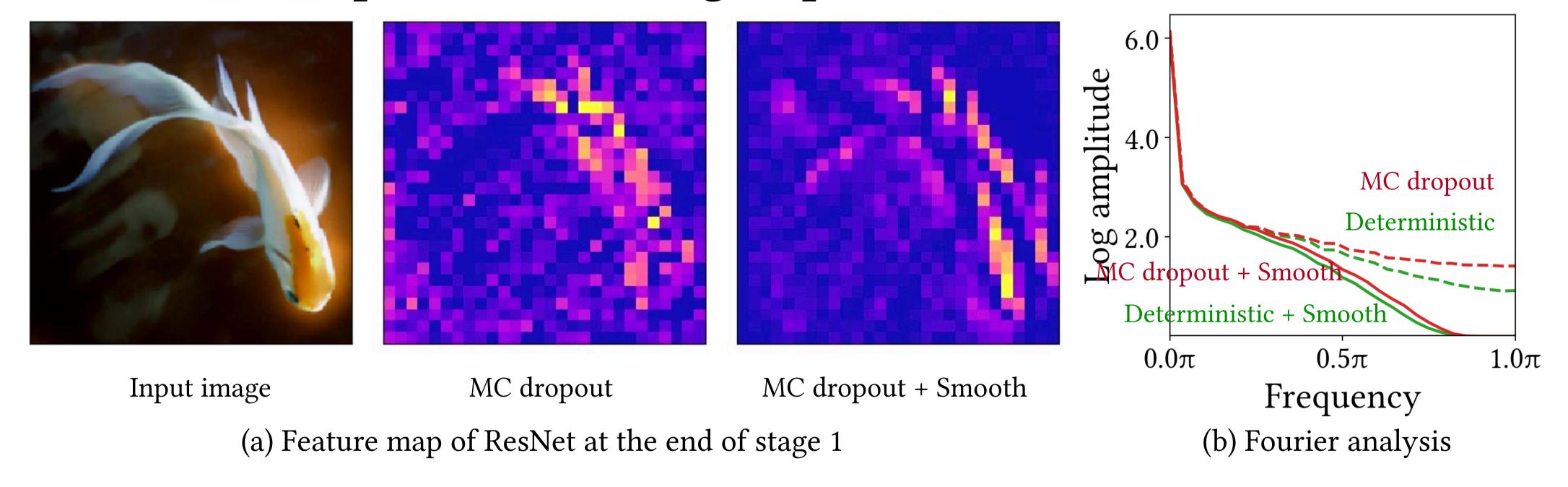


Figure 3. MC dropout adds high-frequency noises. Spatial smoothing **filters high-frequency signals** and stabilizes (denoises) feature maps.

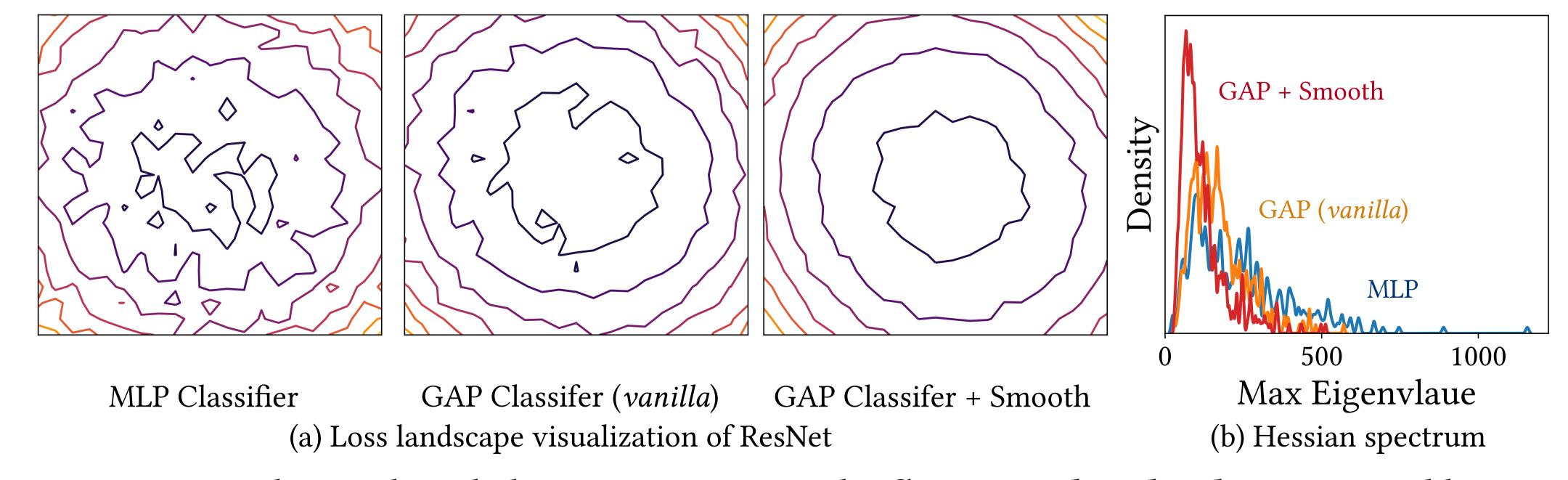


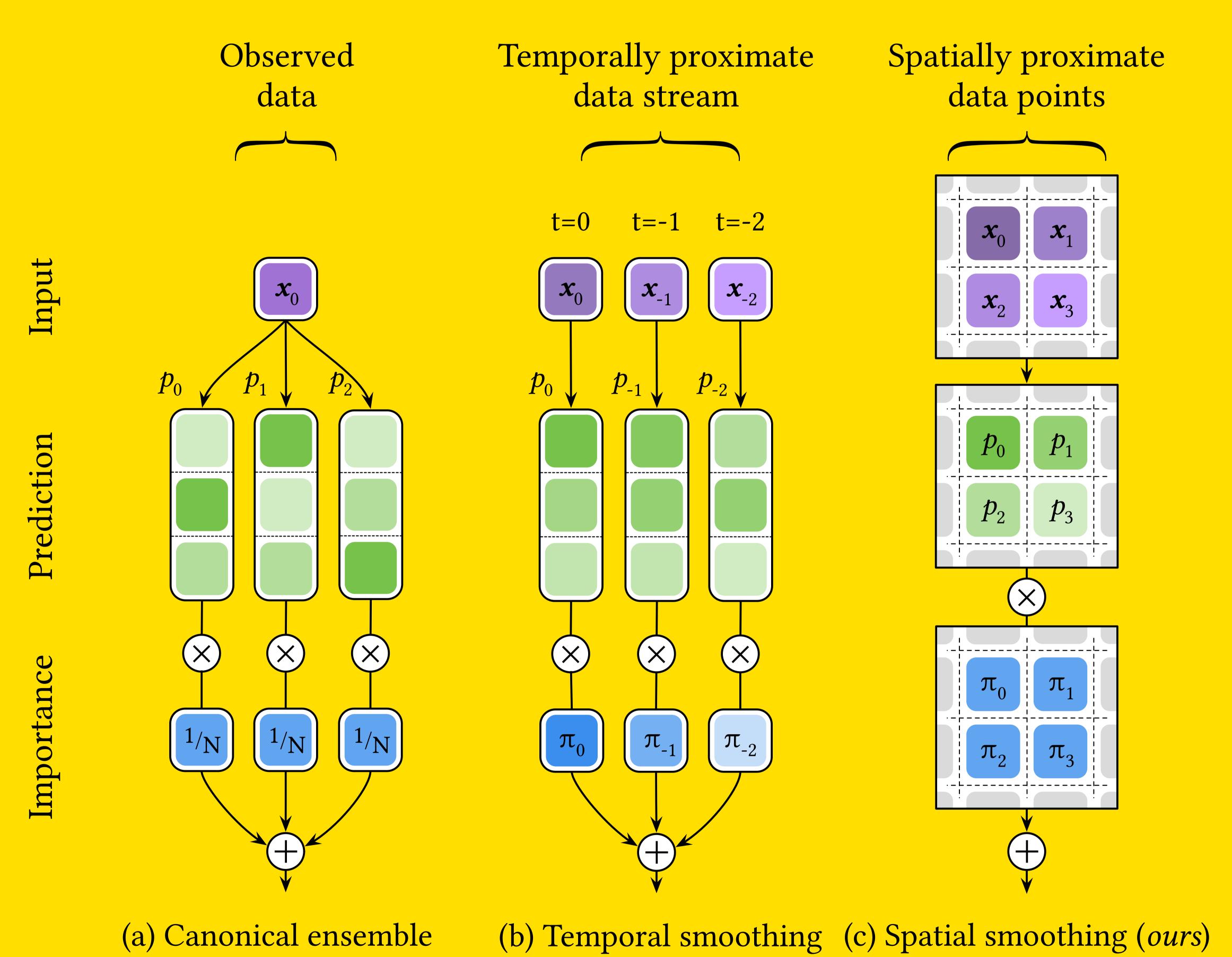
Figure 4. Spatial smoothing helps NN optimization by **flattening loss landscapes**. In addition, GAP is an extreme case of spatial smoothing.

**In summary**, spatial smoothing significantly improves NNs, and has the following properties:

- ① Spatial smoothing flattens loss landscapes.
- ② Spatial smoothing is a low-pass filter.
- ③ Spatial smoothing at the end of a stage plays a key role.



# Exploiting SPATIAL CONSISTENCY is important, and BLURS BEHAVE LIKE ENSEMBLES



# NAMUK PARK, SONGKUK KIM NAVER AI Lab, Yonsei University

# How Do Vision Transformers Work?

We provide an explanation of how *self-attentions* work by addressing them as *a trainable spatial smoothing* of feature maps, because the formulation suggests that self-attentions average feature map values with the positive importance-weights.

# I. Self-Attentions Flatten Loss Landscapes

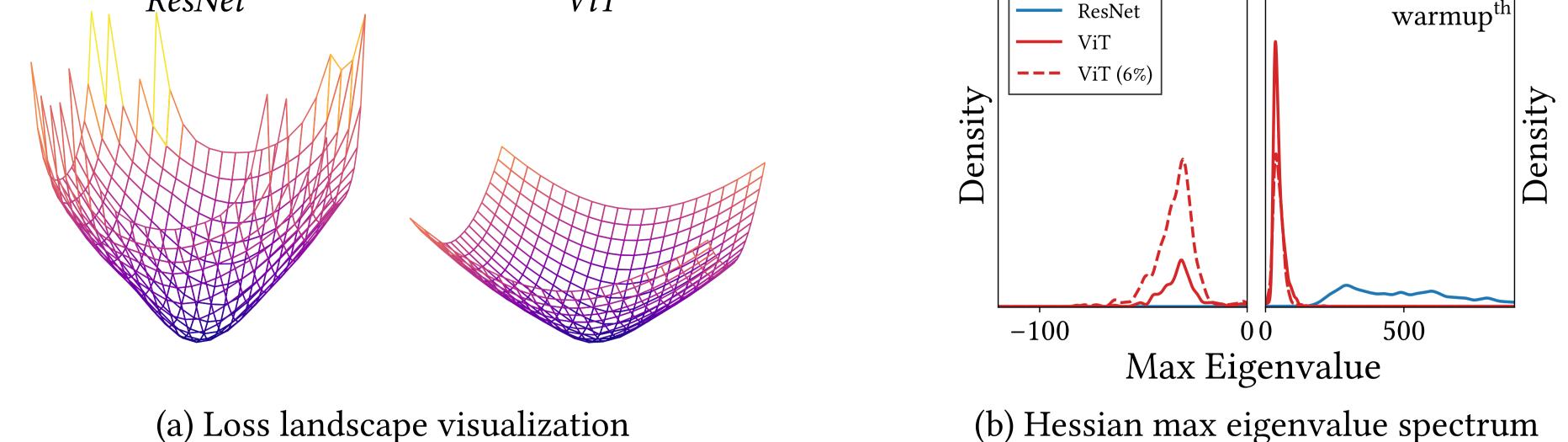


Figure 5. Loss landscape visualization (*Left*) and Hessian max eigenvalue spectrum (*Right*) consistently show that **ViT has a flatter loss than ResNet**.

### II. Self-Attentions Are Low-Pass Filters

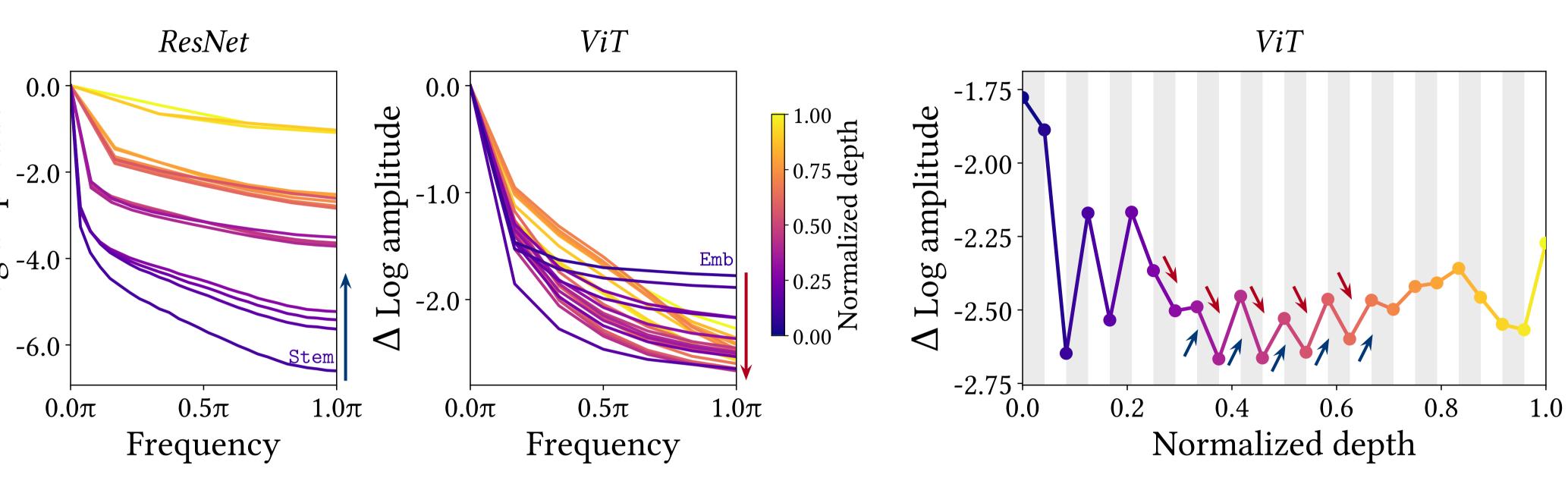


Figure 6. Relative log amplitudes of Fourier transformed feature map show that **ViT tend to reduces high-frequency signals, while ResNet amplify them**. *Right*: In ViT, MSAs (gray area) reduce the high-frequency  $(1.0\pi)$  component, and Conv/MLPs (white area) amplify it.

# III. Self-Attentions at the End of a Stage Play a Key Role in Prediction

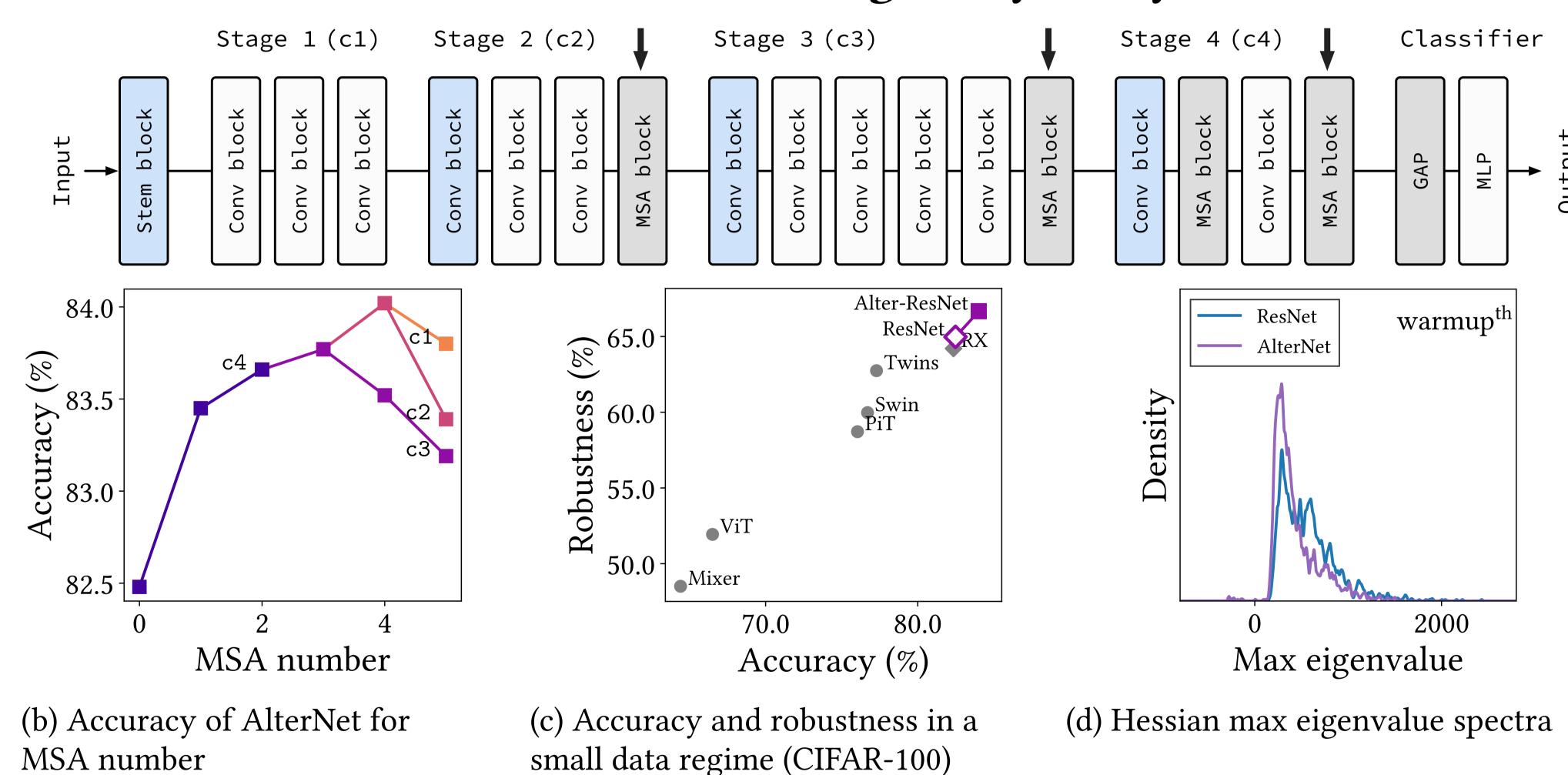
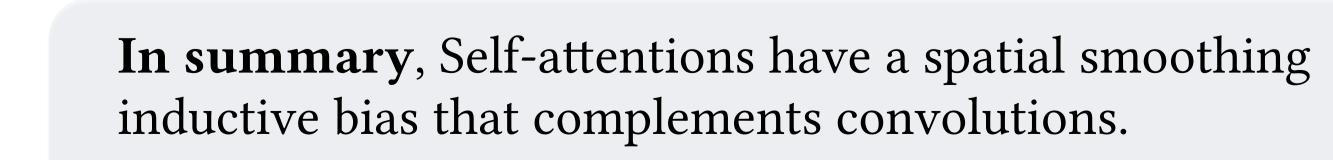


Figure 7. **Self-attentions at the end of a stage (not a model) and Convs at the beginning of a stage significantly improve the performance**. AlterNet, a model in which Conv blocks at the end of a stage are replaced with MSA blocks, outperforms CNNs even on CIFAR.



	Self-Attention	Convolution
Loss Landscape	Flat but non-convex	Convex but sharp
Fourier Analysis	Low-pass filter (shape-biased)	High-pass filter (texture-biased)
<b>Best Practice</b>	The end of a stage	The beginning of a stage

